

Towards Multimodal In-Context Learning for Vision & Language Models

Sivan Doveh^{1,2} Shaked Perek¹ M. Jehanzeb Mirza³ Wei Lin⁵
 Amit Alfassy¹ Assaf Arbelle¹ Shimon Ullman² Leonid Karlinsky⁴

¹IBM Research ²Weizmann Institute of Science

³ICG, TU Graz ⁴MIT-IBM Watson AI Lab ⁵ELLIS Unit, LIT AI Lab,
 Institute for Machine Learning, JKU Linz, Austria

Abstract. State-of-the-art Vision-Language Models (VLMs) ground the vision and the language modality primarily via projecting the vision tokens from the encoder to language-like tokens, which are directly fed to the Large Language Model (LLM) decoder. While these models have shown unprecedented performance in many downstream zero-shot tasks (*e.g.*, image captioning, question answers, *etc.*), still little emphasis has been put on transferring one of the core LLM capability of In-Context Learning (ICL). ICL is the ability of a model to reason about a downstream task with a few examples demonstrations embedded in the prompt. In this work, through extensive evaluations, we find that the state-of-the-art VLMs somewhat lack the ability to follow ICL instructions. In particular, we discover that even models that underwent large-scale mixed modality pre-training and were implicitly guided to make use of interleaved image and text information (intended to consume helpful context from multiple images) under-perform when prompted with few-shot demonstrations (in an ICL way), likely due to their lack of *direct* ICL instruction tuning. To enhance the ICL abilities of the present VLM, we propose a simple yet surprisingly effective multi-turn curriculum-based learning methodology with effective data mixes, leading up to a significant 21.03% (and 11.3% on average) ICL performance boost over the strongest VLM baselines and a variety of ICL benchmarks. Furthermore, we also contribute new benchmarks for ICL evaluation in VLMs and discuss their advantages over the prior art.

1 Introduction

A little more than a year ago, with the release of ChatGPT¹ in late November 2022, Large Language Models (LLMs) made their historical debut showing, for the first time, that an artificial neural network can encompass in its parameters a potentially human-like understanding of language, essentially in all aspects of its distribution complexity including knowledge [19], reasoning [28], context understanding [26] and other core capabilities [31, 64].

¹ <https://chat.openai.com/>

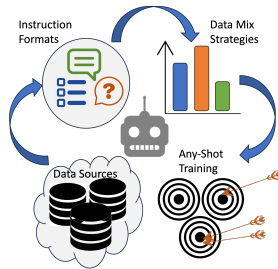


Fig. 1: Multiple data sources are used to generate multi-modal ICL instructions varying the types of ICL tasks and type of semantic concepts shared within each instruction, teaching the VLM to properly correlate information between ICL in-context shots. Our insights on the best training data mix along with our proposed “any-shot” training paradigm enhance the VLM’s ICL abilities.

Besides many other exciting LLM capabilities and discovered emerging properties, perhaps one of the more useful aspects of LLMs is their ‘foundation modeling’ aspect - being able to understand and respond to language input, they are essentially ‘open’ models and are not limited to any set of pre-defined tasks prescribed by the model training. In other words, as part of a typical use - one can explain a task to the model and have it comply (zero-shot inference), or if performance is not sufficiently satisfactory, one can provide a set of *In-Context* demonstrations, illustrating to the model the desired task and input/output format (structure) via a few examples. The latter emerging property of LLMs is commonly referred to as In-Context Learning (ICL) or Few-Shot Learning (FSL) and is extremely useful in situations where more fine-grained control over the downstream task is needed.

Motivated by the great advances of LLMs in language modeling [8, 11, 21, 54, 55], there has been a huge interest in the community towards fusing these LLMs with other modalities. Notably, these LLMs have recently been *fused* with the vision modality [1, 9, 27, 39, 40, 67], Audio [18, 42], Speech [51], Documents [58] *etc.* The resulting models, *e.g.*, the recent Vision and Language Models (VLMs), have demonstrated parallel capabilities to ones demonstrated for language - namely an ability to handle arbitrary downstream tasks in a zero-shot manner, by only explaining the task to the model. While many Vision and Language (VL) benchmarks have been proposed [19, 31, 64], they mostly focus on zero-shot capabilities and tasks defined by common natural language terms, and with few exceptions *e.g.*, [31], do not test the capability of VLMs to leverage paired, image+text, in-context demonstrations for the downstream task definitions. Moreover, as demonstrated by our findings (in Sec. 4), on an exhaustive set of ICL evaluations, some contributed by our work, even the most advanced and recent VLMs [3, 27, 39, 40, 53] still struggle with ICL. This trend seems to be consistent for both the *leading visual encoder to LLM decoder* alignment methods [39, 40], as well as for ‘encoder-free’ techniques that directly processing a mixed stream of multi-modal tokens by a single transformer [3, 27, 53]. Recent VLMs [3, 27, 53]

show some gains on downstream tasks, which can be solved with semantically unrelated in-context demonstrations such as answers to unrelated VQA questions on other images. However, these SOTA VLMs commonly fail when explicitly challenged with tasks (*e.g.*, fine-grained few-shot visual recognition or instance recognition) completely defined by the ICL demonstrations (*e.g.*, an n -way episode of a few-shot visual recognition task). Arguably, ICL is one of the central capabilities of foundation models of all kinds (LLMs and VLMs alike) allowing for fine-grained control over downstream tasks. Hence, such shortcomings of leading VLMs need special attention.

In this work, we conjecture that the ICL performance of modern VLM approaches can be significantly improved by simple modifications to their training strategies, explicitly incorporating *semantically-coherent* ICL tasks into their visual instruction tuning phase. We show that leveraging the *human-assistant* turn-based conversation structure common to visual instruction tuning [39, 40], while extending it to a multi-image conversation, provides a simple and convenient vehicle for multi-shot explicit ICL training. In the multi-turn conversation format, the standard Causal Language Modeling (CLM) objective trains the model to operate in *any-shot* scenario, that is later able to accept any number of in-context demonstration ‘shots’ at inference time. The multi-turn conversation format also includes a ‘zero-shot’ turn (the first turn of the conversation that under CLM object does not have any demonstrations in its processing context) thus providing replay support for zero-shot tasks and avoiding forgetting of the VLM’s core capabilities. Equipped with this adaptation of the visual instruction tuning, we explore and provide insights on the most effective data mixing strategies via forming semantically coherent ICL tasks from multiple available data sources, guided by the requirement to have a common semantic aspect shared by all in-context demonstrations and currently trained query alike. Notably, even the approaches that are trained with mixed, multi-modal, tokenized stream [3, 27, 53] containing multiple images in the same stream, under-perform (as shown in Sec. 4) the above simple training technique proposed by us on top of the more light-weight visual instruction tuning alignment of [40]. As we conjecture, this is likely due to a lack of explicit ICL semantic coherence in their training data design.

To summarize our contributions are as follows: (i) We design a simple and yet surprisingly effective ICL visual instruction tuning strategy that can be easily added to standard visual instruction alignment tuning, significantly enhancing the explicit ICL capabilities of the VLM without forgetting its core zero-shot capabilities; (ii) We analyze and report insights on the most effective data mixes for our proposed ICL instruction tuning; (iii) We offer a set of ICL benchmarks that can assist in testing the ICL abilities of the present-day VLMs and can also act as a standard benchmark for the future.

2 Related Work

We first provide an overview of zero-shot vision-language foundation models and then describe the literature closer to our line of work, *i.e.*, studying in-context learning in the domain of VLMs.

Vision-Language Foundation Models: Recently, VLMs have been adopted as the default choice for *train once and use everywhere* paradigm, and have shown unprecedented performance for many vision-language understanding tasks, *e.g.*, zero-shot classification, visual question-answering (VQA), image captioning, and many more. VLMs can be divided into two families of methods. One family of methods relies on dual-encoders (vision and text encoder) and usually trains the encoders with a contrastive objective by using a large corpus of paired image-text data scraped from the web. Some representatives of this family of methods are CLIP [50] (the first large-scale vision-language model), ALIGN [20], OpenCLIP [52], SigLip [63] and MetaCLIP [61]. Furthermore, some methods have focused their attention on filtering noisy captions (*e.g.*, BLIP [33]), employing textual nearest-neighbors [34] or relying on using geometrically consistent representations [17], and caption augmentations [14, 15] for improving compositional reasoning aspects of these VLMs. In parallel, other methods have employed few-shot supervision [22, 65, 66], and also label-free finetuning [2, 36, 44, 45]. The other group of methods aligns the visual modality with a frozen LLM. BLIP-2 [32] bridges the modality gap between a pre-trained visual encoder and an LLM by using a Querying Transformer. Instruct-BLIP [12] proposes to improve [32] by employing instruction tuning. MiniGPT [67] grounds a vision encoder with a frozen LLM (Vicuna [11]) by only using a trainable linear projection layer between the two. MiniGPT-V2 [9] replaces the LLM with Llama-2 [55] and enhances the performance by also training and finetuning the LLM decoder. Llava [41] also grounds an LLM with a pre-trained visual encoder and also proposes Visual Instruction Tuning, by carefully curating instruction-response pairs, to enhance the performance. The base Llava is further enhanced in Llava-1.5 [37] by careful curation of data and Llava-1.6 [39] also improves the previous version by incorporating some design changes and also modifying the instruction tuning data. Some other works [4, 10, 41, 47, 59] also explore similar ideas. Although these powerful and versatile encoder-decoder models can solve many tasks efficiently, their ability to learn or adapt to new tasks by only seeing a few contextual examples, instead of relying on a huge corpus of training data, is still under-explored. In our work, we take steps towards unlocking the ICL ability of these VLMs by a simple yet effective methodology of carefully curating ICL-specific data and altering the learning framework such that these models can become efficient in-context learners.

In-context Learning for VLMs: In the natural language processing (NLP) literature, the ability of LLMs to solve novel tasks by only consuming a few demonstrations of the downstream task of interest has been formalized as in-context learning (ICL) [43]. For NLP, many different methods have been proposed to

elicit the ICL ability in the powerful autoregressive models for many downstream tasks. Notably [7] popularized few-shot learning for LLMs. Similarly, other methods follow the basic idea of few-shot learning but achieve the goal in different ways, *e.g.*, by breaking down a complex set of instructions in simpler steps [60, 62]. Recently, the vision-language community has also shown considerable interest in ICL. Flamingo [1] showed that ICL can scale up to large-scale vision language models. In particular, they improved downstream tasks like image captioning by only requiring a few examples. Flamingo’s ability was unlocked by their novel fusion of visual information with the textual tokens. Recently, Emu2 [53] showed that ICL for VLMs can be enhanced by scaling up the encoder-decoder models in modern VLMs with auto-regressive training. Similarly, Idefics [27] also scales up the vision (encoder) and language (decoder) to 80 billion and shows effective ICL learning ability. Recent methods have made interesting progress toward obtaining better in-context learning by scaling the model size and specifically training for this purpose [53]. However, we ask the question: can we take an off-the-shelf VLM (like Llava [39]) and convert it to an effective few-shot learner in an ad-hoc fashion? We answer this question in the affirmative and detail our approach in the following sections.

3 Method

In this section, we explain the proposed ICL-instruction-alignment approach. For ease of assimilation, we divide its description into 3 parts as illustrated in Fig. 1. Sec. 3.1 discusses the details of our ICL alignment framework implemented as multi-turn ICL conversations inside the Llava [40] visual-instruction alignment. Sec. 3.2 explains the different ICL instruction task types used in our ICL alignment instruction sets and mixes. Sec. 3.3 discusses the sources of data we used to construct semantically coherent ICL instructions that, as opposed to mixed multi-modal jointly tokenized internet data used in other works [3, 27, 53], is designed to guarantee the existence of a semantic concept shared between each set of ICL demonstrations (shots) and the ICL query. Finally, Sec. 3.4 briefly discusses additional ICL evaluations contributed by us for evaluating ‘open vocabulary few-shot visual recognition’ - one of the most common ICL tasks, somewhat neglected by the previous ICL benchmarks.

3.1 Multi-turn ICL conversations

A common strategy for aligning different pre-trained other-modality encoders (for visual, audio, speech, *etc.*) to an LLM is via multi-modal instruction tuning and associated training curriculum [9, 12, 18, 38–40, 67]. Specifically, we build on the alignment model architecture of [38] which consists of a large-scale pre-trained modality encoder \mathcal{E} , a modality projector \mathcal{P} , and an LLM decoder \mathcal{D} . The training data consists of ‘multi-modal’ conversations between a ‘human’ and a ‘gpt-assistant’, the conversation is typically interleaving these two roles. The human role may contain modality tags (*e.g.*, $\langle image \rangle$ in [38]), yet in practice [38]

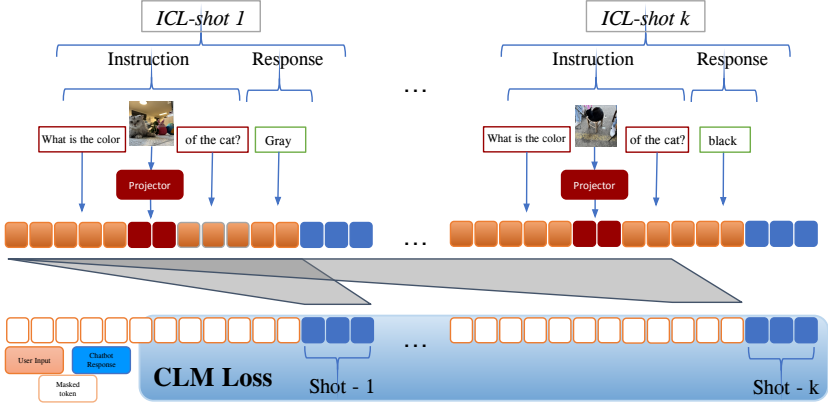


Fig. 2: Causal (left only) attention and formatting the ICL examples as consecutive conversation turns, results in ‘any-shot’ training where the first turn prediction is “zero-shot”, the next turn predicts the response given the context of the first, and so on, resulting in a dynamic “any-shot” context. The grey shades illustrate the context that each turn’s response attends to. As [38] we do completion-only training, masking all but the desired responses (blue) in the target.

only implements this tag once (with a single image per conversation context) added only in the first (human) conversation role text. Therefore, for [38], the training samples are formed as:

$$\mathcal{D}(T_1 \oplus \mathcal{P}(\mathcal{E}(I)) \oplus T_2) \quad (1)$$

where T_1 and T_2 are tokenized parts of the conversation texts that come before and after the $\langle image \rangle$ tag respectively, and \oplus stands for concatenation. Typically, the training curriculum consists of two stages: (i) pre-training - freezing the modality encoder \mathcal{E} and the LLM decoder \mathcal{D} and training the projector \mathcal{P} from scratch; and (ii) fine-tuning - both the projector \mathcal{P} and the LLM decoder \mathcal{D} jointly, keeping \mathcal{E} frozen. In [38] stage (i) training is comprised of short, single-turn (one user, one gpt) conversations, all para-phrasing a request by the user to describe a given image. At the same time, stage (ii) comprises more diverse multi-turn conversations combining different image task instructions and responses, albeit all relative to a single provided context image. In both phases of training, all the user input tokens, *i.e.*, the instruction, and the aligned image tokens, $\mathcal{P}(\mathcal{E}(I))$, are masked and only the gpt-assistant’s responses are used as target labels to train the alignment model with the standard Causal Language Modeling (CLM) objective. Fig. 1 visually describes the process.

We build our ICL visual instruction tuning as a simple and direct extension to [38] alignment tuning. We keep the same architectural components, namely: \mathcal{E} , \mathcal{P} and \mathcal{D} . We keep the pre-training stage (i) intact, yet extend (or replace) the fine-tuning stage (ii) train data with a mix of semantically-coherent ICL instructions. The discussion on the ICL task types and the sources of semantic

coherence is covered in later Sec. 3.2 and Sec. 3.3, while the format of the ICL instructions is as follows:

Human: $S_1^1 \langle \text{image} \rangle S_1^2$ GPT: R_1
 Human: $S_2^1 \langle \text{image} \rangle S_2^2$ GPT: R_2
 Human: $S_3^1 \langle \text{image} \rangle S_3^2$ GPT: $R_3 \dots$

accompanied with a corresponding ordered list of images $[I_1, I_2, I_3, \dots]$. Here each $\{S_j^1, S_j^2, I_j, R_j\}$ comprises an ICL ‘shot’, that is a single in-context demonstration example. Such multi-turn ICL instruction format required only minor modifications to the official [38] code². The only modification was to add multi-image input support in training as opposed to the original [38] training phases, which only used single-image conversations. Interestingly, in combination with input masking of human turns, the CLM objective, and the CLM attention - allowing tokens to attend only to their left, our simple ICL conversation format becomes an ‘any-shot trainer’. Indeed, the first shot has no other shots in its attended context and therefore serves as a ‘zero-shot’ instruction replay - effectively reminding the aligned model that it needs to continue supporting this mode of operation, as aligned with [38]. Then, each shot number i , observes the shots $1, \dots, (i-1)$ in its attended context and hence trains the model to support $(i-1)$ -shot ICL instructions.

3.2 ICL instruction task types

We set each ICL instruction shot $\{S_j^1, S_j^2, I_j, R_j\}$ to one of the following instruction-response formats: (a) Open QA - the texts S_j^1 and S_j^2 form an open question, where image I_j tokens are embedded in the semantically valid position prescribed by the question language context. The desired response R_j is formulated as natural language text (commonly a single sentence); (b) Multiple-choice QA - S_j^1 is empty and S_j^2 contains a question followed by several answer options marked by A, B, C, ..., R_j contains a single letter correct answer choice; (c) Captioning - S_j^1 contains a para-phrased in different ways ‘describe the image’ request, S_j^2 is empty, and R_j contains the image description. For all ICL instruction task types, the shots in the same instruction have a semantic coherence in the form of a semantic concept shared across all the shots in the ICL instruction (e.g., all questions ask about a certain type of object attributes such as color, or all captions share a common style and/or intent - e.g. describe locations of objects on the image). The coherence is achieved via careful data curation explained in Sec. 3.3. Intuitively, in a single ICL instruction, all the shots are of the *same* instruction task type (a) Open question answering (QA), (b) multiple choice (MC), or (c) Captioning (Cap), respectively. More details and examples of each ICL instruction type are provided in Supplementary.

² <https://github.com/haotian-liu/LLaVA>

3.3 Data sources for ICL instruction mixes

As we show in our ablation studies in Sec. 5, correct mixing of ICL instruction types, as well as sources of semantic coherence, is crucial for maximizing ICL instruction alignment performance. In other words, the ICL curriculum is an important component of our approach. To induce semantic coherence between shots of each ICL instruction, that is, the presence of some semantic concept shared between all shots in each any-shot ICL conversation explained in Sec. 3.1, we collect our ICL instructions from the following data sources: SEED benchmark [31] partitions 1-5 (Scene Understanding, Instance Identity, Instance Attributes, Instance Location, Instances Counting) and VL-Checklist [64] 13 partitions (*e.g.*, Color, Material, Action, Object, Positional Relation, Action Relation, Action Attribute). In all cases, we use the source dataset partition to sample the k -shots of each ICL instruction from the same partition. As described in Sec. 3.2 all shots of the ICL instruction are then structured in the same format according to one of the 3 formats described. As extensively analyzed in Sec. 5, the data mixes of instruction formats and sources of semantic coherence (the aforementioned dataset partitions) have a crucial importance on downstream ICL performance and generalization. We contribute insights into these empirical observations on the optimal data curriculums (resulting from significant investments in computing) with hopes of inspiring exciting future research direction of enhancing ICL in VLMs. Additionally, we found that aside from the first (zero-shot) turn replay, adding the original [38] fine-tune data portion to the training curriculum for providing additional replay support is useful.

3.4 ICL benchmarks

With few exceptions (*e.g.*, task 23 in Seed-2 [30]), most of ICL evaluation benchmarks so far were constructed ad-hoc from largely unrelated k -shot episodes randomly sampled from either a VQA dataset (VQAr2, VQA, OKVQA, TextVQA, VizWiz, *etc.*) or a visual dialog dataset [13], as evaluated by [27, 53]. While such evaluations demonstrate improvements with adding more shots, the lack of semantic relations between the in-context demonstration shots and the query makes it unclear whether most of the improvement does not come from, *e.g.*, matching the desired output format. Additionally, practical ICL use cases are often parallel to the now classical problem of few-shot learning, particularly few-shot visual recognition. Equipped with these insights, we formulate these additional ICL benchmarks for VLMs: (i) Fine-grained few-shot visual recognition ICL benchmarks derived from popular few-shot tasks: Stanford Dogs [23], CUBS [57], FOOD-101 [5], Stanford Cars [24], Flowers [46]; (ii) ‘unseen’ Seed tasks 6-8 reformulated into ICL episodes as described above (but not used for training); (iii) put aside validation portion of ‘seen’ Seed task 5 (instance counting) as the hardest of the training Seed tasks requiring all the other capabilities of the training Seed tasks 1-5 to execute; (iv) put aside validation portion of ‘seen’ VL-checklist partitions. We detail the exact statistics of all these datasets and splits in the Supplementary.

4 Results

In Sec. 4.1, we first provide an overview of the different datasets used in our work, followed by a brief description of the baselines and an explanation of the implementation details. We discuss our main findings in Sec. 4.2, uncovering the strong potential for our proposed ICL instruction tuning for decoder VLMs, leading to over 11% average absolute improvement over the strongest baseline. We provide an ablation study in Sec. 5 digging deeper into the different aspects of our method, primarily focusing on the properties of the ICL instruction tuning data mixes that lead to the aforementioned improvements. Additionally, we highlight the scaling potential of our approach by increasing the ICL instruction tuning data size. Finally, we ablate the importance of replaying non-ICL data instruction tuning data, concluding that (a) due to non-ICL instructions replay, our model avoids forgetting the core capabilities of the base model [38] as measured by the standard and extensive MME [16] benchmark; (b) the replay visual instruction tuning data generally improves our model’s ICL performance.

4.1 Evaluation Settings

Datasets: We briefly list the datasets used to form our ICL-instruction tuning mixes and/or for evaluation of ICL or other capabilities of the resulting model and baselines.

1. **SEED-Bench-2** [30] - SEED-Bench-2 [30] is an extended version of SEED-Bench [31] that features a total of 27 evaluation dimensions. Our use of SEED-Bench-2 is two-fold. We use 100% of data of tasks 1-4 (scene and ‘Instance’ tasks) to form ICL instructions of multiple-choice type (SEED-Bench format). We use 90% of data of tasks 5 (‘Instances Counting’, arguably hardest of 1-5 tasks) to form multiple-choice ICL instructions and its 10% to form 2-shot multiple-choice ICL ‘Instances Counting’ evaluation. Additionally, we use 100% of tasks 6-8 (relation, interaction, reasoning) for 2-shot multiple-choice ICL evaluation of ‘Unseen SEED-tasks’. Additionally, we evaluate directly on SEED-Bench-2 task 23 - ‘In-context captioning’, the only SEED ICL task. We provide more details on ICL with SEED-Bench-2 in the Supplementary.
2. **VL-checklist** [64] - is a benchmark constructed from Visual Genome [25], SWiG [49], VAW [48], and HAKE [35]. Overall, VL-checklist has 13 such partitions according to a shared semantic aspect of MM evaluation: attribute, relation, or object, each further subdivided. We split VL-checklist into two non-overlapping subsets, using 70% for ICL instruction tuning and 30% for ICL 2-shot testing. For VL-checklist we form both multiple-choice and open QA ICL task types for all partitions and, additionally, ICL captioning task types for 6 partitions: Color, State, Material, Size, and Action attribute/relation. More details and examples are provided in the Supplementary.
3. **LLaVA visual instruction tuning dataset** [38] - dataset of 655K visual instructions constructed by the LLaVA team [40] as part of their visual

instruction tuning works [38–40]. Built from a combination of COCO, GQA, OCR-VQA, TextVQA, and VisualGenome data. We do not use this data for ICL training, only for general visual instructions replay intended for preserving the general capabilities of the base [38] VLM which we build upon.

4. **Stanford Dogs** [23] - is a dataset of 120 dog breeds with 150 images per class, 12,000 images for training and 8,580 for testing.
5. **CUB** [56] - contains 200 images of different types of bird species with 5,994 samples for training and 5,794 for testing. The annotations include several attributes and localization.
6. **Flowers** [46] - consists of 102 images of flower categories common in the UK. Each class contains between 40 and 258 images. Its test set has 6149 images.
7. **Food-101** [6] - is a data set of 101 food categories, with 100K images. Each class contains 250 manually reviewed test images and 750 training images. The full test set includes 25250 images.
8. **Stanford Cars** [24] - consists of 196 classes of cars with a 50/50 division to train and test. Categories are typically at the level of Make, Model, and Year. Test set contains 8041 samples.

Datasets 4-8 are fine-grained few-shot datasets that are completely unseen during the training of our model and are used for evaluating our ICL instruction tuning generalization testing only. We test them in a 2-way / 1-shot mode, meaning that for each testing episode (ICL task instance), in addition to the test image, we sample a random image from the same class as the test image and a random image from a different class than the test image. We then provide them to the model with their respective labels, and finally, we provide the model with the test image as a query asking the model to choose among the two class possibilities. We use only the test partitions of these datasets and generate an episode for each test sample from the test sets. We generate an average of 10763 episodes from each dataset.

Baselines: We compare with the following state-of-the-art baselines.

1. IDEFICS 9B [27] - is an open-source reproduction of Flamingo [1], trained on a large corpus of publicly available datasets, and shows strong in-context learning abilities. The architectural details are kept similar to [1]; however, the total number of model parameters differ.
2. OpenFlamingo [3] - also proposes to provide an open-source alternative to the original Flamingo [1], and sticks to the design proposed by [1], while training on different data than the original Flamingo.
3. EMU2 [53] - is a generative multimodal model capable of generating both images and texts. It has good ICL abilities and has been shown to achieve better performance than IDEFICS-80B and Flamingo-9B on multiple tasks.
4. LLaVA-1.5 13B [38] - builds upon the framework proposed by the original LLaVA [40] but makes some design changes, like using a more expressive

projection between the vision and the language tokens and also trains on different instruction-tuning data.

5. LLaVA-next (1.6) 13B [39] - improves LLaVA-1.5 by altering the instruction tuning data and the visual processing pipeline.

Metrics: All experiments, unless stated otherwise, were measured using the accuracy of the models’ generation abilities, requiring an exact match of the GT and the predicted model. The accuracy is calculated as the percentage of exact-match responses. The response may be a single character, as in the multi-choice or few-shot scenarios, or a full string in the captioning and open-question tests. An exception to the exact-match criteria is the SEED-23 task, where the output is measured using the perplexity value, selecting the most likely character out of the four choices (the choices are provided as part of task 23; one of the choices gives the desired caption that should be selected given the context of 2 image examples and their respective semantically related captions).

Implementation Details: We build upon the [38] codebase and use its default training parameters when tuning our model. We extended the codebase to allow accepting a list of encoded images with each visual instruction during training (the original code supported only a single image tag and image input per visual instruction conversation). Following this extension, a training conversation is allowed to contain a number of <image> tags matching the length of the associated images list. These <image> tags are replaced according to the order of the images list. We use Vicuna-1.5-13B [11] LLM backbone and use the recommended by [38] 2K context length, which supports up to 3-images conversations for our any-shot ICL instruction tuning (each image takes about 500 tokens of the context). Larger models and more extensive investment in computing would allow much larger context length and support for much longer any-shot sequences, which we believe would improve the performance further. We used a single 8x A100 80GB Nvidia GPU node for all of our training runs.

4.2 Main findings

In this section, we discuss our findings from a comprehensive set of multi-modal ICL evaluation experiments carried out on a collection of held-out validation splits of our ICL instruction tuning mix datasets (SEED-Bench-2 task 5 (Instance counting), VL-checklist partitions), completely unseen ICL tasks evaluating our method generalization (SEED-Bench-2 tasks 6-8, 5 few-shot fine-grained visual recognition datasets), as well as native ICL captioning evaluation of SEED-Bench-2: task 23 (currently the only semantically coherent multi-modal ICL evaluation task).

Few-shot visual recognition evaluations: We evaluate our model and the base-lines (discussed in Sec. 4.1) on the collection of 5 standard (fine-grained) few-shot datasets (listed in Sec. 4.1) using 2-way / 1-shot episodes of multiple-choice ICL

Model	Food	Cars	Dogs	CUB	Flowers	AVG
IDEFICS 9B	65.94	77.30	50.93	62.00	55.29	59.30
OpenFlamingo 9B	52.3	57.43	50.47	51.2	48.78	52.04
EMU2 37B	59.92	55.42	50.27	53.56	52.76	52.47
Llava-1.5 13B	87.19	57.30	33.00	58.24	58.60	60.77
Otter	33.15	22.19	33.8	24.62	60.01	34.75
Llava-1.6 13B	89.15	84.70	72.39	67.90	65.58	72.96
Ours 13B	97.44	96.51	79.85	78.67	76.51	85.79
(gain)	(+8.29)	(+11.8)	(+7.46)	(+10.77)	(+10.93)	(+12.83)

Table 1: Comparison across our proposed fine-grained few-shot ICL tasks. The last row highlights the gains compared to Llava-1.6 13B, which is the leading baseline.

task type. All these datasets are unseen during training, and test our method’s multi-modal ICL generalization ability and the baselines. The results of this evaluation are presented in Tab. 1. Clearly, our method was able to leverage the proposed multi-turn successfully conversation-based any-shot ICL tuning and the ICL data mix to generalize well to these unseen fine-grained few-shot visual recognition tasks over a diverse set of visual categories from the 5 datasets, including food, vehicles, animals and plants. On average, our method successfully improves by over 12% over the top-performing baseline ([39]), notably improving over IDEFICS 9B [27], Otter [29] and EMU2 [53] by over 25% on average despite of their strong, mixed-modality, pre-training which includes multiple images in their respective training contexts. Intuitively, we attribute these large gains to a likely lack of explicit semantic coherence in the baselines’ training data. Interestingly and expectedly, our semantically coherent ICL pertaining generalizes well to these tasks, showing good potential for further improvement by introducing our proposed ICL instruction tuning to all LLM-decoder-based VLMs and highlighting the importance of fixing the ICL performance in those otherwise very strong models.

Additional ICL evaluations: We present additional ICL evaluations in Tab. 2. These are comprised of a mix of ICL validation set tasks (SEED-Bench-2 Instance Counting, multiple-choice / QA / captioning ICL on VL-checklist) and unseen (during training) ICL generalization tasks (SEED-Bench-2 tasks 6-8 multiple-choice ICL and in-context captioning task 23). Similar findings to those observed for the fine-grained few-shot ICL experiments in Tab. 1 are seen. Again, our method improves by over 10 points beyond the top-performing baseline, further highlighting the importance of our contributions and the need to incorporate the proposed semantically coherent ICL training into LLM-decoder-based VLMs. Notably, our method improves over 4 points in SEED-Bench-2 task 23 (in-context captioning) - the only semantically coherent ICL benchmark task to the best of our knowledge that was not contributed by us, further confirming our intuition above. We note, on the low results of the OpenFlamingo model, that in many cases, the model returns empty values, which adds up as mistakes. On average, across all datasets, we see an average boost of 11%.

Model	SEED	23	Unseen	Ins Count	MC	VL	QA	VL	Cap	VL	AVG
IDEFICS 9B	45.00	27.00	28.00	63.77	10.00	2.37	29.36				
OpenFlamingo 9B	25.00	24.80	18.25	45.00	0.77	0.41	19.04				
EMU2 37B	-	44.15	30.9	67.60	12.00	2.25	31.38				
Llava-1.5 13B	44.17	61.00	60.00	33.00	1.56	5.62	34.26				
Otter	34.7	27.09	34.27	65.79	12.75	2.27	29.47				
Llava-1.6 13B	40.8	59.00	49.00	88.00	19.48	2.27	43.09				
Ours 13B	49.16	62.20	65.30	95.37	24.50	23.40	53.32				
(gain)	(+4.16)	(+3.20)	(+16.30)	(+7.37)	(+5.02)	(+21.03)	(+10.23)				

Table 2: Evaluating ICL tasks built from SEED and VL-Checklist. MC=multiple choice, QA=open QA, Cap=ICL captioning. Ins Count=our SEED-2 task 5 val split. Unseen=our ICL test based on SEED-2 tasks 6-8. SEED 23=SEED-2 task 23. EMU2 did not report on SEED and their public code does not support perplexity inference.

5 Ablations

In this section, we explore the different aspects and design choices that contribute to the success of our approach of ICL instruction tuning. In Sec. 5.1 we explore the effects of different choices and mixes of *task types* and *shared semantic concepts* of the ICL instructions on the performance of the model. In Sec. 5.2 we discuss the scaling properties of our model with respect to adding more ICL data. In Sec. 5.3 we explore how well our additional ICL instruction tuning preserves the base capabilities of the [38] model we start from. In particular, we show that replaying non-ICL instruction data (from [38]) helps preserve the model capabilities (as measured through the MME [16] metrics), and also, surprisingly, generally benefits performance. Finally, in Sec. 5.4 we analyze the model’s ability to effectively leverage in-context (visual + text) information during inference. We measure the model’s performance by ranging the number of shots between 2 and 0, showing that as desired, our ICL-instructions-tuned model strongly benefits from the addition of more shots (in-context visual and text information).

5.1 ICL Data mixes ablation

We evaluate multiple different ICL instruction data mixing strategies and the effect of including the base [38] fine-tuning data into our mix. These experiments are summarized in Tab. 3. We will refer to the different mixes via their mix ID. As we can see from comparing mixes 1 and 2, and our best mix 5 with mix 6, including [38] data has a positive effect on the performance of the ICL tasks (and also preserves the aligned model capabilities). We, therefore, include it in all the rest of the data mixing recipes. Next, we discuss the aspects of ICL instruction type and shared semantic concepts in the ICL instructions. We use the average across all evaluations from Tab. 1 and Tab. 2 as our average accuracy measure. Here, we derive some important insights that significantly boost our ICL instruction tuning performance (as described below), yet of course, with more investment in compute performance could be significantly improved

Mix ID	LLaVA data	Attributes	Relations	Categories	Instances	Open Questions	Multiple Choice	Captioning	AVG
1		0.00%	0.00%	0.00%	100.00%	0.00%	100.00%	0.00%	56.56
2	✓	0.00%	0.00%	0.00%	100.00%	0.00%	100.00%	0.00%	65.41
3	✓	66.67%	0.00%	0.00%	33.33%	50.00%	50.00%	0.00%	66.35
4	✓	75.00%	12.50%	0.00%	12.50%	65.00%	5.00%	30.00%	68.44
5	✓	45.45%	15.15%	36.36%	3.04%	39.40%	42.42%	18.18%	69.33
6	✓	46.87%	15.62%	37.05%	0.00%	40.625%	40.625%	18.75%	68.10
7		45.45%	15.15%	36.36%	3.04%	39.40%	42.42%	18.18%	67.84

Table 3: Shared semantic concept (between ICL shots in the same any-shot instruction) mixing ablation (left). Instruction format mixes effect on model performance ablation (right). ‘LLaVA data’ indicates if tuning data from [38] was used in the mix.

	Cognition	Total Perception	Total
Llava1.5	295.36	1531	
Ours	301.43	1520	

Table 4: MME scores of baseline and our model.

	MC	QA	Cap
Two-shot	95.37	24.5	23.4
One-shot	87.9	12.9	0
Zero-shot	87.4	12.	0

Table 5: Varying number of Shots on VL tasks.

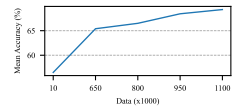


Fig. 3: Mean Accuracy (%) while scaling ICL instruction data.

further. As our experiments already show great promise for our proposed ICL instruction tuning, we leave the investigation to future work.

Mixing ICL instruction formats: Here we examine the importance of the ICL instruction format in ICL tuning examples formation (Tab. 3 right side). Namely, we measure the relative importance of open questions vs multiple-choice vs captioning ICL instructions. Comparing mixes 2 - 5, starting from multiple-choice only ICL instructions, adding open questions ICL poses a significant benefit to performance. Additionally, we find that including ICL captioning task helps performance, and the best mix that is behind our best-performing model reported in Tab. 1 and Tab. 2, is reported under mix ID 5.

Shared semantic concepts within ICL instructions: Part of our ICL instruction design is having a shared semantic concept (category, instance, relation, attribute, etc.) in all the shots of a single ICL instruction. This design teaches the model to be sensitive and be able to leverage the shared semantic information between shots. Here we analyze which types of semantic concepts composition in the mix is most beneficial. According to our findings (Tab. 3 left side), shared attribute concepts seem to be most beneficial, as long as it doesn’t completely dominate the data. The second most important concept type is ‘object categories’ (e.g., having cat notion in all shots), while interestingly ‘object instances’ (e.g., asking about particular cat instances in all shots) seems less helpful to the best mix. Yet removing instances data completely hurts performance (mix

6). Relations shared concepts are represented in the third place, yet comparing mixes with and without them, we see they provide a significant boost.

5.2 Data scaling

Here we evaluate the scaling potential of our ICL instruction tuning approach. Fig. 3 presents our model average performance with scaling the ICL instruction tuning data. As can be seen, more ICL data consistently improves average performance over ICL tasks, ending with a positive gradient indicating the potential for additional improvement with further scaling.

5.3 Preserving the base model abilities

One of the important questions is, with the addition of the ICL instructions - are we able to preserve the base capabilities of the aligned VLM? In particular, the [38] we are starting from. We answer this question in the affirmative in Tab. 4, comparing our ICL instruction tuned model MME [16] scores with the base [38] model. We attribute the preservation of the model’s base capabilities to the importance of replaying the [38] data, which was also found effective for boosting the ICL tasks performance in Sec. 5.1.

5.4 Leveraging the shots information

Finally, in Tab. 5 we test our model performance on a subset of our ICL evaluation tasks changing the number of in-context shots between 0 and 2 (the maximum allowed by the set 2K context length of [38]). As we see, our model positively improves with the addition of more shots, as expected.

6 Summary & Conclusions

In this work, we analyzed and proposed some ways to improve the ICL tasks performance of Vision and Language Models. Our proposed approach leverages carefully designed ICL instructions and their respective data mixes, as well as the proposed any-shot training paradigm resulting in a model able to take advantage of in-context examples to better perform on a variety of tasks, such as multiple choice Q&A, instance counting, captioning etc. Our work includes extensive comparisons to strong baselines. We also propose and evaluate few-shot visual recognition posed as ICL on multiple fine-grained datasets. Simple as it is, our approach shows significant and consistent gains across all evaluations, suggesting that future VLMs can significantly benefit from the proposed ICL instruction tuning (with shared semantic concepts) as well as from the any-shot training paradigm. We believe our ideas are orthogonal to our current implementation and can be easily re-used for many other models. Our work suggests exciting future work directions including exploring ICL instruction tuning with longer LLM context enabling longer any-shot sequences, additional exploration of ICL instruction data mixes, additional ICL task types, and more ideas on possible shared ICL semantics.

References

1. Alayrac, J.B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., Ring, R., Rutherford, E., Cabi, S., Han, T., Gong, Z., Samangooei, S., Monteiro, M., Menick, J., Borgeaud, S., Brock, A., Nematzadeh, A., Sharifzadeh, S., Binkowski, M., Barreira, R., Vinyals, O., Zisserman, A., Simonyan, K.: Flamingo: a Visual Language Model for Few-Shot Learning. arXiv:2204.14198 (2022) [2](#), [5](#), [10](#)
2. Alfassy, A., Arbelle, A., Halimi, O., Harary, S., Herzig, R., Schwartz, E., Panda, R., Dolfi, M., Auer, C., Saenko, K., Staar, P.J., Feris, R., Karlinsky, L.: Feta: Towards specializing foundation models for expert task applications (2022) [4](#)
3. Awadalla, A., Gao, I., Gardner, J., Hessel, J., Hanafy, Y., Zhu, W., Marathe, K., Bitton, Y., Gadre, S., Sagawa, S., Jitsev, J., Kornblith, S., Koh, P.W., Ilharco, G., Wortsman, M., Schmidt, L.: Openflamingo: An open-source framework for training large autoregressive vision-language models. arXiv preprint arXiv:2308.01390 (2023) [2](#), [3](#), [5](#), [10](#)
4. Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., Zhou, J.: Qwen-vl: A frontier large vision-language model with versatile abilities. arXiv preprint arXiv:2308.12966 (2023) [4](#)
5. Bossard, L., Guillaumin, M., Van Gool, L.: Food-101 – mining discriminative components with random forests. In: European Conference on Computer Vision (2014) [8](#), [22](#)
6. Bossard, L., Guillaumin, M., Van Gool, L.: Food-101 – mining discriminative components with random forests. In: European Conference on Computer Vision (2014) [10](#)
7. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Nee-lakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. Advances in neural information processing systems **33**, 1877–1901 (2020) [5](#)
8. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Nee-lakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D.: Language Models are Few-Shot Learners. arXiv:2005.14165 (2020) [2](#)
9. Chen, J., Li, D.Z.X.S.X., Zhang, Z.L.P., Xiong, R.K.V.C.Y., Elhoseiny, M.: MINIGPT-V2: Large Language Model As a Unified Interface for Vision-language Multi-task Learning. arXiv preprint arXiv:2310.09478 (2023) [2](#), [4](#), [5](#)
10. Chen, K., Zhang, Z., Zeng, W., Zhang, R., Zhu, F., Zhao, R.: Shikra: Unleashing multimodal llm’s referential dialogue magic. arXiv preprint arXiv:2306.15195 (2023) [4](#)
11. Chiang, W.L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J.E., Stoica, I., Xing, E.P.: Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality (March 2023), <https://lmsys.org/blog/2023-03-30-vicuna/> [2](#), [4](#), [11](#)
12. Dai, W., Li, J., Li, D., Tiong, A., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: InstructBLIP: Towards general-purpose vision-language models with instruction tuning. In: Thirty-seventh Conference on Neural Information Processing Systems (2023) [4](#), [5](#)

13. Das, A., Kottur, S., Gupta, K., Singh, A., Yadav, D., Moura, J.M., Parikh, D., Batra, D.: Visual dialog. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 326–335 (2017) [8](#)
14. Doveh, et al.: Teaching structured vision & language concepts to vision & language models. In: CVPR (2023) [4](#)
15. Doveh, S., Arbellet, A., Harary, S., Alfassy, A., Herzig, R., Kim, D., Giryes, R., Feris, R., Panda, R., Ullman, S., et al.: Dense and aligned captions (dac) promote compositional reasoning in vl models. arXiv preprint arXiv:2305.19595 (2023) [4](#)
16. Fu, C., Chen, P., Shen, Y., Qin, Y., Zhang, M., Lin, X., Yang, J., Zheng, X., Li, K., Sun, X., et al.: Mme: A comprehensive evaluation benchmark for multimodal large language models. arXiv preprint arXiv:2306.13394 (2023) [9](#), [13](#), [15](#)
17. Goel, S., Bansal, H., Bhatia, S., Rossi, R.A., Vinay, V., Grover, A.: CyCLIP: Cyclic Contrastive Language-Image Pretraining. arXiv:2205.14459 (2022) [4](#)
18. Gong, Y., Luo, H., Liu, A.H., Karlinsky, L., Glass, J.: Listen, think, and understand. arXiv preprint arXiv:2305.10790 (2023) [2](#), [5](#)
19. Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., Steinhardt, J.: Measuring massive multitask language understanding. Proceedings of the International Conference on Learning Representations (ICLR) (2021) [1](#), [2](#)
20. Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q.V., Sung, Y., Li, Z., Duerig, T.: Scaling Up Visual and Vision-Language Representation Learning With Noisy Text Supervision. In: Proc. ICML (2021) [4](#)
21. Jiang, A.Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D.S., Casas, D.d.l., Hanna, E.B., Bressand, F., et al.: Mixtral of experts. arXiv preprint arXiv:2401.04088 (2024) [2](#)
22. Khattak, M.U., Rasheed, H., Maaz, M., Khan, S., Khan, F.S.: Maple: Multi-modal Prompt Learning. In: IEEE Conference on Computer Vision and Pattern Recognition (2023) [4](#)
23. Khosla, A., Jayadevaprakash, N., Yao, B., Fei-Fei, L.: Novel dataset for fine-grained image categorization. In: First Workshop on Fine-Grained Visual Categorization, IEEE Conference on Computer Vision and Pattern Recognition. Colorado Springs, CO (June 2011) [8](#), [10](#), [22](#)
24. Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3d object representations for fine-grained categorization. In: Proceedings of the IEEE international conference on computer vision workshops. pp. 554–561 (2013) [8](#), [10](#), [22](#)
25. Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Li, L.J., Shamma, D.A., et al.: Visual genome: Connecting language and vision using crowdsourced dense image annotations. International journal of computer vision **123**, 32–73 (2017) [9](#)
26. Lampinen, A.K., Dasgupta, I., Chan, S.C., Matthewson, K., Tessler, M.H., Creswell, A., McClelland, J.L., Wang, J.X., Hill, F.: Can language models learn from explanations in context? arXiv preprint arXiv:2204.02329 (2022) [1](#)
27. Laurençon, H., Saulnier, L., Tronchon, L., Bekman, S., Singh, A., Lozhkov, A., Wang, T., Karamcheti, S., Rush, A.M., Kiela, D., Cord, M., Sanh, V.: Obelics: An open web-scale filtered dataset of interleaved image-text documents (2023) [2](#), [3](#), [5](#), [8](#), [10](#), [12](#)
28. Lewkowycz, A., Andreassen, A., Dohan, D., Dyer, E., Michalewski, H., Ramasesh, V., Slone, A., Anil, C., Schlag, I., Gutman-Solo, T., et al.: Solving quantitative reasoning problems with language models. Advances in Neural Information Processing Systems **35**, 3843–3857 (2022) [1](#)
29. Li, B., et al.: Mimic-it: Multi-modal in-context instruction tuning. arXiv (2023) [12](#)

30. Li, B., Ge, Y., Ge, Y., Wang, G., Wang, R., Zhang, R., Shan, Y.: Seed-bench-2: Benchmarking multimodal large language models (2023) [8](#), [9](#)
31. Li, B., Wang, R., Wang, G., Ge, Y., Ge, Y., Shan, Y.: Seed-bench: Benchmarking multimodal llms with generative comprehension. arXiv preprint arXiv:2307.16125 (2023) [1](#), [2](#), [8](#), [9](#), [21](#)
32. Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597 (2023) [4](#)
33. Li, J., Li, D., Xiong, C., Hoi, S.: BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. arXiv:2201.12086 (2022) [4](#)
34. Li, Y., Liang, F., Zhao, L., Cui, Y., Ouyang, W., Shao, J., Yu, F., Yan, J.: Supervision Exists Everywhere: A Data Efficient Contrastive Language-Image Pre-training Paradigm. arXiv:2110.05208 (2021) [4](#)
35. Li, Y.L., Liu, X., Wu, X., Li, Y., Qiu, Z., Xu, L., Xu, Y., Fang, H.S., Lu, C.: Hake: A knowledge engine foundation for human activity understanding. TPAMI (2023) [9](#)
36. Lin, W., Karlinsky, L., Shvetsova, N., Possegger, H., Kozinski, M., Panda, R., Feris, R., Kuehne, H., Bischof, H.: Match, expand and improve: Unsupervised finetuning for zero-shot action recognition with language knowledge. In: ICCV (2023) [4](#)
37. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning. arXiv preprint arXiv:2310.03744 (2023) [4](#)
38. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning (2023) [5](#), [6](#), [7](#), [8](#), [9](#), [10](#), [11](#), [13](#), [14](#), [15](#)
39. Liu, H., Li, C., Li, Y., Lee, Y.J.: Llava-next (llava 1.6) (2023), <https://llava-vl.github.io/blog/2024-01-30-llava-next/> [2](#), [3](#), [4](#), [5](#), [10](#), [11](#), [12](#)
40. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. In: NeurIPS (2023) [2](#), [3](#), [5](#), [9](#), [10](#)
41. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. arXiv preprint arXiv:2304.08485 (2023) [4](#)
42. Lyu, C., Wu, M., Wang, L., Huang, X., Liu, B., Du, Z., Shi, S., Tu, Z.: Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. arXiv preprint arXiv:2306.09093 (2023) [2](#)
43. Min, S., Lewis, M., Zettlemoyer, L., Hajishirzi, H.: Metaicl: Learning to learn in context. arXiv preprint arXiv:2110.15943 (2021) [4](#)
44. Mirza, M.J., Karlinsky, L., Lin, W., Possegger, H., Feris, R., Bischof, H.: TAP: Targeted Prompting for Task Adaptive Generation of Textual Training Instances for Visual Classification. arXiv preprint arXiv:2309.06809 (2023) [4](#)
45. Mirza, M.J., Karlinsky, L., Lin, W., Possegger, H., Kozinski, M., Feris, R., Bischof, H.: LaFTer: Label-free tuning of zero-shot classifier using language and unlabeled image collections. In: NeurIPS (2023) [4](#)
46. Nilsback, M.E., Zisserman, A.: Automated flower classification over a large number of classes. In: 2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing. pp. 722–729 (2008). <https://doi.org/10.1109/ICVGIP.2008.47> [8](#), [10](#), [22](#)
47. Peng, Z., Wang, W., Dong, L., Hao, Y., Huang, S., Ma, S., Wei, F.: Kosmos-2: Grounding multimodal large language models to the world. arXiv preprint arXiv:2306.14824 (2023) [4](#)
48. Pham, K., Kaffe, K., Lin, Z., Ding, Z., Cohen, S., Tran, Q., Shrivastava, A.: Learning to predict visual attributes in the wild. In: Proceedings of the IEEE/CVF Con-

- ference on Computer Vision and Pattern Recognition (CVPR). pp. 13018–13028 (June 2021) [9](#)
49. Pratt, S., Yatskar, M., Weihs, L., Farhadi, A., Kembhavi, A.: Grounded situation recognition. *ArXiv abs/2003.12058* (2020) [9](#)
 50. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning Transferable Visual Models from Natural Language Supervision. In: *Proc. ICML* (2021) [4](#)
 51. Radford, A., Kim, J.W., Xu, T., Brockman, G., McLeavey, C., Sutskever, I.: Robust speech recognition via large-scale weak supervision. In: *International Conference on Machine Learning*. pp. 28492–28518. PMLR (2023) [2](#)
 52. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C.W., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., Schramowski, P., Kundurthy, S.R., Crowson, K., Schmidt, L., Kaczmarczyk, R., Jitsev, J.: LAION-5b: An open large-scale dataset for training next generation image-text models. In: *NeurIPS* (2022) [4](#)
 53. Sun, Q., Cui, Y., Zhang, X., Zhang, F., Yu, Q., Luo, Z., Wang, Y., Rao, Y., Liu, J., Huang, T., et al.: Generative multimodal models are in-context learners. *arXiv preprint arXiv:2312.13286* (2023) [2](#), [3](#), [5](#), [8](#), [10](#), [12](#)
 54. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and Efficient Foundation Language Models. *arXiv:2302.13971* (2023) [2](#)
 55. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al.: Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023) [2](#), [4](#)
 56. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: Cub-200-2011 (caltech-ucsd birds-200-2011). Tech. Rep. CNS-TR-2011-001, California Institute of Technology (2011) [10](#)
 57. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011 dataset. Tech. rep., California Institute of Technology (2011) [8](#), [22](#)
 58. Wang, D., Raman, N., Sibue, M., Ma, Z., Babkin, P., Kaur, S., Pei, Y., Nourbakhsh, A., Liu, X.: Doclm: A layout-aware generative language model for multimodal document understanding. *arXiv preprint arXiv:2401.00908* (2023) [2](#)
 59. Wang, W., Chen, Z., Chen, X., Wu, J., Zhu, X., Zeng, G., Luo, P., Lu, T., Zhou, J., Qiao, Y., et al.: Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175* (2023) [4](#)
 60. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q.V., Zhou, D., et al.: Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems* **35**, 24824–24837 (2022) [5](#)
 61. Xu, H., Xie, S., Tan, X.E., Huang, P.Y., Howes, R., Sharma, V., Li, S.W., Ghosh, G., Zettlemoyer, L., Feichtenhofer, C.: Demystifying clip data. In: *Proc. ICLR* (2023) [4](#)
 62. Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., Narasimhan, K.: Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems* **36** (2024) [5](#)
 63. Zhai, X., Mustafa, B., Kolesnikov, A., Beyer, L.: Sigmoid loss for language image pre-training. *arXiv preprint arXiv:2303.15343* (2023) [4](#)
 64. Zhao, T., Zhang, T., Zhu, M., Shen, H., Lee, K., Lu, X., Yin, J.: Vl-checklist: Evaluating pre-trained vision-language models with objects, attributes and relations. *arXiv preprint arXiv:2207.00221* (2022) [1](#), [2](#), [8](#), [9](#), [21](#)
 65. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Conditional Prompt Learning for Vision-Language Models. In: *Proc. CVPR* (2022) [4](#)

- 66. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to Prompt for Vision-Language Models. IJCV (2022) [4](#)
- 67. Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M.: Minigpt-4: Enhancing Vision-language Understanding with Advanced Large Language Models. arXiv preprint arXiv:2304.10592 (2023) [2](#), [4](#), [5](#)

Appendix

We first provide details about the training data and the ICL instruction tasks (Section A), then list the detailed ICL results per task (Section B). Later, provide details about the ICL instruction task types during the test phase (Section C) and finally conclude with qualitative visualizations.

A ICL instruction task types in training

We construct our training data from two public datasets, namely SEED benchmark [31] and VL Checklist [64]. Next, we provide details regarding how the training data is formalized.

A.1 SEED benchmark

Data from tasks 1-5 from the SEED benchmark is used for creating training in-context learning instructions that share a common semantic concept within each ICL instruction, see Figs. 4 and 5 for qualitative examples. Furthermore, we use task 1-4 data only for ICL training purposes while keeping a portion of task 5 data for ICL evaluation on Instance Counting (IC). Additionally, we generate ICL tasks from SEED benchmark tasks 6-8 for testing only (to test generalization to unseen ICL tasks) as explained in Appendix C.2. In Tab. 11 we provide the detailed statistics of the SEED-bench train data used in our train split.

A.2 VL-Checklist

Data from tasks 0-12 from the VL-Checklist benchmark is used for creating training in-context learning instructions with a shared common semantic concept within each ICL instruction, see figures Figs. 6 and 7 for examples. We split data in each task to non-overlapping train and test portions to test the resulting VLM on the ICL capability for each semantic concept (task) of the VL-Checklist (as explained in Appendix C.3). In Tab. 12 we provide the detailed statistics of the VL-Checklist train data used in our train split.

B ICL benchmarks- results per task

B.1 VL Checklist

In Tabs. 6 to 8, we provide results per task on the test splits of the VL-Checklist ICL tasks. We separately provide results for different ICL task types: QA-question answering in Tab. 6, MC - multiple choice questions in Tab. 7, and Cap - captioning in Tab. 8.

B.2 SEED benchmark

In Tab. 9, we provide results per task on the test partitions of the SEED benchmark-based ICL tasks we created.

C ICL instruction task types in test

We show in Figs. 12 to 14 examples from each dataset and task used in the test set.

C.1 Few-shot benchmarks

We create fine-grained, few-shot visual recognition ICL benchmarks from the following fine-grained classification datasets: Stanford Dogs [23], CUBS [57], FOOD-101 [5], Stanford Cars [24], Flowers [46].

C.2 SEED benchmark tests examples

For ICL testing on SEED benchmark, we use a test split for SEED benchmark task 5 (Instance Counting) that we created (‘seen ICL task’ - observed during training that contained a train partition of task 5), as well as the entire data of SEED benchmark tasks 6-8 (used for testing only, checking generalization of our ICL instruction tuning approach to ‘unseen ICL tasks’ - unobserved during training). Fig. 14 provides examples of the unseen ICL tasks from SEED benchmark - 6 7 8. We are showing these examples already formatted as ICL tasks. In Tab. 10 we provide the detailed statistics of the SEED-bench test data used in our test split.

C.3 VL-Checklist test examples

For our study, we utilized the dataset provided by the VL-Checklist paper. We meticulously divided it into training and test sets through a random selection process, ensuring that there was no overlap of images between the two sets. We provide examples in the train-set VL checklist Figs. 6 and 7. For this reason, we will not provide examples from the test set, as it would be redundant. In Tab. 12 we provide the detailed statistics of the VL-Checklist test data used in our test split.

Model	material	size	action	color	state	RelA	RelS	L	S	M	center	margin	mid
Ours	0.35	0.36	0.23	0.46	0.27	0.12	0.26	0.13	0.17	0.23	0.13	0.19	0.29

Table 6: Accuracy (%) per Question answering task in VL Checklist. RelA - relation action, RelS-relation spatial, L-Large Object, S-Small Object, M-Medium Object.

Model	material	size	action	color	state	RelA	RelS	L	S	M	center	margin	mid
Ours	0.99	0.76	0.96	0.94	0.95	0.98	0.98	0.96	0.98	0.97	0.97	0.99	0.97

Table 7: Accuracy (%) per Multiple Choice task in VL Checklist. RelA - relation action, RelS-relation spatial, L-Large Object, S-Small Object, M-Medium Object.

Model	material	size	action	color	state	RelA
Ours	0.21	0.19	0.36	0.18	0.21	0.26

Table 8: Accuracy (%) per captioning task in VL Checklist.

Model	IC	SR	II	VR
Ours	0.65	0.53	0.74	0.76

Table 9: Accuracy (%) per task in SEED Bench. IC-Instance Counting, SR- Spatial Relation, II-Instance Interaction, VR- Visual Reasoning.

Split	IC	SR	II	VR
Test	251	657	97	331

Table 10: Amount of items in test split on SEED Bench (we use part of task 5 and all of the data of tasks 6-8 for testing). IC-Instance Counting, SR- Spatial Relation, II-Instance Interaction, VR- Visual Reasoning. The ‘Test’ row in the table provides the amount of test images in each task.

Split	SU	II	IA	IL	IC
Train	3158	1831	4649	978	2196

Table 11: Amount of items in train split on SEED Bench (we use all data of tasks 1-4 and part of the data of task 5 for training). SU-Scene Understanding, II-Instance Identity, IA-Instance Attribute, IL-Instance Location, IC-Instance Counting.

Model	material	size	action	color	state	RelA	RelS	L	S	M	center	margin	mid
Train	12000	12000	4437	12000	5292	12000	900	12000	12000	12000	12000	12000	12000
Test	1431	1802	494	7317	588	12000	100	7915	2805	2999	7929	2485	7901

Table 12: Amount of samples in train and test splits of VL-Checklist tasks.

HUMAN: What kind of city is the background of the image?
A. Rural B. Historical C. Modern D. Coastal
Answer with the option's letter from the given choices directly.

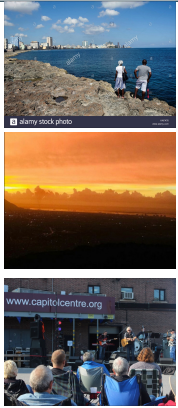
GPT: D

HUMAN: What time of day is captured in the image?
A. Morning B. Noon C. Night D. Sunset or sunrise
Answer with the option's letter from the given choices directly.

GPT: D

HUMAN: What is the main event taking place in the image?
A. A group of people playing soccer in a park B. A crowd of people watching a live music performance C. A group of people attending a political rally D. A street performer playing music while people pass by
Answer with the option's letter from the given choices directly.

GPT:



www.capitolcentre.org

Answer: B

(a) SEED benchmark task 1: Scene Understanding

HUMAN: In the image, what is next to the wooden bench?
A. A metal fence B. A green fence C. A wooden fence D. A brick wall
Answer with the option's letter from the given choices directly.


GPT: B

HUMAN: Is there any bike visible in the image?
A. Yes B. Blurred info C. Can't say D. No
Answer with the option's letter from the given choices directly.

GPT: D

HUMAN: What is the primary ingredient that can be seen next to the small white bowl of tomato sauce?
A. Tomato sauce B. Tomatoes C. Spinach D. Basil
Answer with the option's letter from the given choices directly.

GPT:



Answer: B

(b) SEED benchmark task 2: Instance Identity

HUMAN: What color is the building in the background?
A. Red B. Gray C. Brick D. White
Answer with the option's letter from the given choices directly.


GPT: D

HUMAN: What is the color of the leather furniture in the image?
A. Gray B. Brown C. Red D. Orange
Answer with the option's letter from the given choices directly.

GPT: B

HUMAN: What is the predominant color of the sand on the beach?
A. White B. Brown C. Black D. Gray
Answer with the option's letter from the given choices directly.

GPT:



Answer: B

(c) SEED benchmark task 3: Instance Attribute

Fig. 4: SEED benchmark training examples tasks 1-3

HUMAN: Where is the spoon in relation to the bowl of food?
A. Under the bowl B. Next to the bowl C. Above the bowl D. Inside the bowl
Answer with the option's letter from the given choices directly.


GPT: D

HUMAN: What is the position of the keyboard in the image?
A. Behind the computer monitor B. In front of the man's face C. On the left side of the desk D. On the right side of the desk
Answer with the option's letter from the given choices directly.

GPT: C

HUMAN: What is the position of the red car in the image?
A. Centered in the middle of the image B. In the foreground, left side of the image C. In the background, right side of the image D. Cannot be determined
Answer with the option's letter from the given choices directly.

GPT:



Answer: A

(a) SEED benchmark task 4: Instance Location

HUMAN: How many objects are visible in the living room?
A. 3 B. 5 C. 4 D. 6
Answer with the option's letter from the given choices directly.

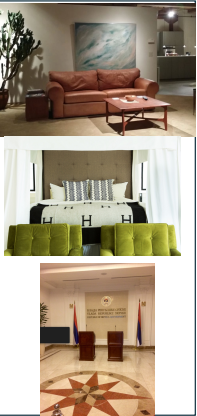
GPT: D

HUMAN: How many pillows are on the bed in the image?
A. 4 B. 3 C. 2 D. 5
B. Answer with the option's letter from the given choices directly.

GPT: A

HUMAN: How many podiums are in the image?
A. One B. Two C. Three D. Four
Answer with the option's letter from the given choices directly.


GPT:



Answer: B

(b) SEED benchmark task 5: Instance Counting

Fig. 5: SEED benchmark training examples tasks 4-5




HUMAN: Provide a one-sentence caption for the provided image
GPT: A kitchen with wooden floors and wooden chairs

HUMAN: Provide a one-sentence caption for the provided image
GPT: A computer desk with two monitors and a keyboard

HUMAN: Provide a one-sentence caption for the provided image
GPT:

Answer: A bicycle with a black frame and orange accents."

(a) VL Checklist Caption task 0: Material



HUMAN: Provide a one-sentence caption for the provided image
GPT: A small boat is on the water

HUMAN: Provide a one-sentence caption for the provided image
GPT: A large green bus is parked at a bus stop.

HUMAN: Provide a one-sentence caption for the provided image
GPT:

Answer: A small bicycle is sitting on top of a bookshelf.

(b) VL Checklist Caption task 1: Size



HUMAN: Provide a one-sentence caption for the provided image
GPT: A woman in a bikini holding a frisbee.




HUMAN: Provide a one-sentence caption for the provided image
GPT: A man with a black beard and mustache is eating a hot dog.

HUMAN: Provide a one-sentence caption for the provided image
GPT:

Answer: A group of giraffes are standing in a forest

(c) VL Checklist Caption task 2: Action

Fig. 6: VL Checklist captioning training examples tasks 0-2






HUMAN: Provide a one-sentence caption for the provided image
GPT: A young boy is pulling a red suitcase.

HUMAN: Provide a one-sentence caption for the provided image
GPT: A large white boat with a blue stripe.

HUMAN: Provide a one-sentence caption for the provided image
GPT:

Answer: A man wearing a grey jacket and red shoes is skiing down a snowy hill.

(a) VL Checklist Caption task 3: Color



HUMAN: Provide a one-sentence caption for the provided image
GPT: A bathroom stall with graffiti on the wall

HUMAN: Provide a one-sentence caption for the provided image
GPT: A baby is laying on a bed with a white comforter

HUMAN: Provide a one-sentence caption for the provided image
GPT:

Answer: A refrigerator is open and full of food and drinks

(b) VL Checklist Caption task 4: State



HUMAN: Provide a one-sentence caption for the provided image
GPT: A man is surfing on a wave in the ocean


HUMAN: Provide a one-sentence caption for the provided image
GPT: A man brushes his teeth while looking at the camera

HUMAN: Provide a one-sentence caption for the provided image
GPT:

Answer: A boy in a yellow shirt is about to kick a soccer ball

(c) VL Checklist Caption task 5: Relative Action

Fig. 7: VL Checklist captioning training examples tasks 3-5




HUMAN: What kind of floor is it? A. paper floor B. wooden floor
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What material is the boat made of? A. porcelain boat B. wooden boat
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What material is used for the printer? A. brass printer B. plastic printer
Answer with the option's letter from the given choices directly
GPT:

Answer: B

(a) VL Checklist Multiple Choice task 0: Material




HUMAN: What is the size of the rock? A. small rock B. large rock
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What is the size of the food? A. large food B. small food
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: What is the size of the tree? A. small tree B. large tree
Answer with the option's letter from the given choices directly
GPT:

Answer: B

(b) VL Checklist Multiple Choice task 1: Size



HUMAN: What are the people doing? A. walking pedestrians B. playing pedestrians
Answer with the option's letter from the given choices directly
GPT: A

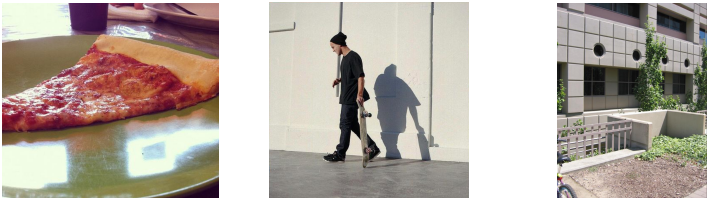
HUMAN: What is the boy doing? A. swimming boy B. skateboarding boy
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What is the activity of the person? A. walking person B. swinging person
Answer with the option's letter from the given choices directly
GPT:

Answer: A

(c) VL Checklist Multiple Choice task 2: action

Fig. 8: VL Checklist Multiple Choice training examples tasks 0-2




HUMAN: What is the color of the plate? A. white plate B. turquoise plate
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: What is the color of the board? A. gray board B. light red board
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: What is the color of the leaves? A. green leaves B. purple leaves
Answer with the option's letter from the given choices directly
GPT:

Answer: A

(a) VL Checklist Multiple Choice task 3: Color




HUMAN: What is the state of the fabric? A. clean fabric B. dusty fabric
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: what is the state of the metal? A. cracked metal B. full metal
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: Are the cans empty? A. full cans B. empty cans
Answer with the option's letter from the given choices directly
GPT:

Answer: A

(b) VL Checklist Multiple Choice task 4: State



HUMAN: What is the person doing with the tie? A. person wear tie B. person assemble tie
Answer with the option's letter from the given choices directly
GPT: A




HUMAN: Is the person holding something? A. person cut surfboard B. person hold surfboard
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What action is person doing to the sports ball? A. person kick sports ball B. person dry sports ball
Answer with the option's letter from the given choices directly
GPT:

Answer: A

(c) VL Checklist Multiple Choice task 5: Relative Action

Fig. 9: VL Checklist Multiple Choice training examples tasks 3-5




HUMAN: Do you see any tracks in the snow? A. tracks out of snow B. tracks **IN** snow
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: Where is the picture in relation to the box? A. picture in front of box B. picture behind box
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: Is there a shadow above or below the boy? A. shadow under boy B. shadow **ON** boy
Answer with the option's letter from the given choices directly
GPT:

Answer: B

Fig. 10: VL Checklist Multiple Choice task 6: Relative Spatial




HUMAN: What is the material of the chair? Answer the question using a single word or phrase
GPT: iron

HUMAN: What is the frame made of? Answer the question using a single word or phrase.
GPT: metal

HUMAN: What material is used for the printer? Answer the question using a single word or phrase
GPT:

Answer: Brass

(a) VL Checklist Question Answering task 0: Material



HUMAN: Is there food on the plate? Answer the question using a single word or phrase.
GPT: no

HUMAN: Is the bench empty or full? Answer the question using a single word or phrase.
GPT: empty

HUMAN: Are the cans empty? Answer the question using a single word or phrase.
GPT:

Answer: No, they are full




(b) VL Checklist Question Answering task 4: State

Fig. 11: VL Checklist Question Answering training examples tasks 0,4

HUMAN: What is the type of the bird in the image?
A.Black footed Albatross B. Least Tern
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: What is the type of the bird in the image?
A.Black footed Albatross B. Least Tern
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What is the type of the bird in the image?
A.Black footed Albatross B. Least Tern
Answer with the option's letter from the given choices directly
GPT:






Answer: A

(a) Few shot CUB

HUMAN: What is the model of the car in the image?
A. Ford Freestar Minivan 2007 B. BMW ActiveHybrid 5 Sedan 2012
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: What is the model of the car in the image?
A. Ford Freestar Minivan 2007 B. BMW ActiveHybrid 5 Sedan 2012
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What is the model of the car in the image?
A. Ford Freestar Minivan 2007 B. BMW ActiveHybrid 5 Sedan 2012
Answer with the option's letter from the given choices directly
GPT:






Answer: B

(b) Few shot Stanford Cars

HUMAN: What is the type of the flower in the image?
A.rose B. pink primrose
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What is the type of the flower in the image?
A.rose B. pink primrose
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: What is the type of the flower in the image?
A.rose B. pink primrose
Answer with the option's letter from the given choices directly
GPT:



Answer: B

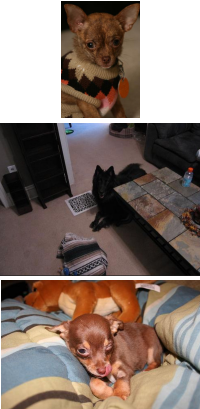
(c) Few shot Flowes

Fig. 12: Few shot datasets

HUMAN: What is the breed of the dog in the image?
A. Chihuahua B. groenendael
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: What is the breed of the dog in the image?
A. Chihuahua B. groenendael
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What is the breed of the dog in the image?
A. Chihuahua B. groenendael
Answer with the option's letter from the given choices directly
GPT:




Answer: A

(a) Few shot Stanford Dogs

HUMAN: What is the type of the food in the image?
A. filet mignon B. churros
Answer with the option's letter from the given choices directly
GPT: B

HUMAN: What is the type of the food in the image?
A. filet mignon B. churros
Answer with the option's letter from the given choices directly
GPT: A

HUMAN: What is the type of the food in the image?
A. filet mignon B. churros
Answer with the option's letter from the given choices directly
GPT:



Answer: B

(b) Few shot Food101

Fig. 13: few shot datasets

HUMAN: What is the position of the saxophones in the band in the image?
A. In the front B. In the middle C. Can't tell from the image D. In the back
Answer with the option's letter from the given choices directly.

GPT: C


HUMAN: What is the relative position of the drums to the person playing them?
A. The drums are in front of the person B. The drums are to the right of the person C. The drums are to the left of the person D. The drums are behind the person
Answer with the option's letter from the given choices directly.

GPT: A

HUMAN: What is the position of the man relative to the camera in the photo?
A. In front of the camera B. Beside the camera C. Behind the camera D. On top of the camera
Answer with the option's letter from the given choices directly.

GPT:

Answer: C



(a) SEED benchmark task 6: Spatial Relation

HUMAN: What is the relative position of the woman and the violin in the foreground?
A. The woman is standing next to the violin B. The woman is sitting on the violin C. The violin is on a stand behind the woman D. The woman is holding the violin
Answer with the option's letter from the given choices directly.

GPT: D

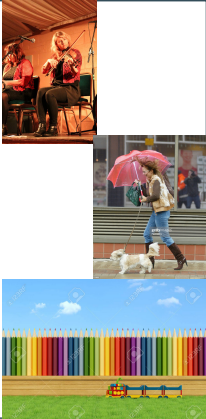
HUMAN: What is the relation between the woman and the dog in the image?
A. The woman is walking her own dog. B. The woman is playing fetch with the dog. C. The dog is a stray and the woman is trying to catch it. D. The woman is training the dog
Answer with the option's letter from the given choices directly.

GPT: A

HUMAN: What is the relation between the train and the colored pencils?
A. The train is passing by a pencil fence B. The train is running on a railroad track made of colored pencils C. The train is pulling a long row of colored pencils D. The train is pushing a long row of colored pencils
Answer with the option's letter from the given choices directly.

GPT:

Answer: A



(b) SEED benchmark task 7: Instance Interactions

HUMAN: If a person is cooking on the stove, which direction should they face?
A. Toward the window B. Toward the floor C. Toward the cabinets D. Toward the sink
Answer with the option's letter from the given choices directly.

GPT: D

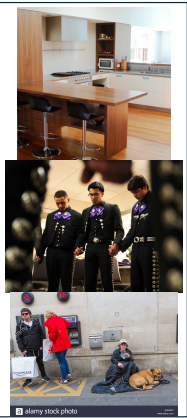
HUMAN: How many men are in the image, and where are they located?
A. Five men in Mexican outfits standing next to each other B. Four men in suits standing in a line, with one touching a chain held by a group of men in Mexican outfits C. Six men in Mexican outfits holding hands D. Three men in Mexican outfits standing far away from each other
Answer with the option's letter from the given choices directly.

GPT: C

HUMAN: What can be inferred about the group of people sitting on the street?
A. They are homeless B. They are street performers C. They are waiting for a parade D. They are tourists
Answer with the option's letter from the given choices directly.

GPT:

Answer: D



(c) SEED benchmark task 8: Visual Reasoning

Fig. 14: SEED benchmark Test examples tasks 6-8